

本体论与信息检索

廖明宏

(哈尔滨工业大学计算机科学与工程系, 哈尔滨150001)

摘 要: 试图对本体论做形式化的描述, 并在此基础上讨论本体论的一个应用: 基于本体论的信息检索方法, 它克服了目前基于关键词检索所造成的信息冗余和丢失的不足。其查询方法更符合人类的思维习惯, 查询结果也更合理、可用。

关键词: 本体论; 信息检索; 企业知识管理

Ontology and Information Retrieval

Liao Minghong

(Department of Computer Science and Engineering, Haerbin Institute of Technology, Haerbin 150001)

【Abstract】 This paper tries to describe the ontology in a formal way and introduce one of its application: an information retrieval method based on ontology. This method overcomes the insufficiency of information overloading and missing caused by the keyword-based retrieval methods. So it more conforms to the custom of human thought, and the results become more rational and useful. This method is verified by an experimental protosystem.

【Key words】 Ontology; Information retrieval; Enterprise knowledge management

本体论(或称实体论)这个概念在计算机科学中扮演着越来越重要的角色^[1]。然而, 到目前为止, 在计算机界还很难为本体论下一个确切的定义, 但斯坦福大学的Gruber给出的定义得到许多同行的认可, 即本体论是对概念化的精确描述^[2]。本体论的最终目标是精确地表示那些隐含(或不明确的)信息, 使得它们可被软件系统重用和共享。

随着因特网的普及与推广, 人们越来越依赖于万维网从事科研及商务活动。因而在网上检索信息成了人们最关心的话题之一。常听到有人抱怨, 利用现有检索工具来查询某一信息, 得到的结果是一堆信息垃圾, 很少有他们想要的东西, 原因在于目前的信息检索工具主要是基于关键词或内容分类目录进行查找的, 很少做进一步的智能化处理。本文借鉴本体论的基本思想, 提出一种基于本体论的信息检索方法。这种方法使查询结果更能满足用户的要求。

1 概念化与本体论

本体论是对概念化的精确描述, 然而本体论与概念化仍存在本质的区别。在人工智能教科书中, 概念化被定义为一个结构 $\langle D, R \rangle$, 其中 D 表示一个域, R 为 D 上的关系。 R 是一般数学意义上的关系, 故称为外延函数(Extensional relations)。然而, 这里所关心的是这些关系本身的含义, 称为内涵关系(Intensional relations)或概念关系(Conceptual relations)。如果说外延关系是定义在某个域上的, 那么内涵关系则是定义在某个域空间上。

定义1: 域空间是一个结构 $\langle D, W \rangle$, 其中 D 表示一个域, W 表示域中元素相关状态的集合(或称可能的世界)。

定义2: n 元概念关系: 域空间 $\langle D, W \rangle$ 上定义的 n 元概念关系是一个从 W 到 D 上所有 n 元关系集合的全函数, 即: W 。

定义3: 概念化: 域 D 的概念化是一个三元组 $C = \langle D, W, R_c \rangle$, 其中, R_c 表示域空间 $\langle D, W \rangle$ 上的概念关系集合。

定义4: 语言的内涵解释是一个结构 $\langle C, IC \rangle$, 其中 $C = \langle D, W, R_c \rangle$ 是一个概念化, 而函数 J 将 D 中的元素赋上 V 中的常量符号, R_c 元素赋上 V 中的谓词符号。

-56-

语言的内涵解释也称为语言的本体论的承诺(Ontological Commitment)。若 $K = \langle C, IC \rangle$ 表示语言 L 的本体论的承诺, 我们说 L 通过 K 承诺 C , 而 C 是 K 的基本概念化。概念化与本体论的联系与区别在于给定一个带有本体论的承诺 K 的语言 L , L 的本体论是一个公理集合, 该公理集合使得它的模型集合尽可能近似于 L 关于 K 的原有模型。事实上这样的公理集合是不易找到的, 因此我们说一个语言 L 的本体论近似于一个概念化 C , 如果存在一个本体论的承诺 $K = \langle C, IC \rangle$, 使得 L 关于 K 的原有模型被包含于本体论的模型之中。

2 基于本体论的信息检索方法

目前的信息检索方法主要是基于关键词或目录分类的, 其查询结果往往产生大量毫无相关的信息, 同时又可能丢失重要的信息。由于本体论刻画了事物之间的内在联系, 借助于本体论, 可以使检索的信息更能满足用户的需求。下面以查询企业雇员能力为例, 介绍一种启发式的信息检索方法。

在现代化企业管理中人们越来越清楚地认识到知识, 尤其是雇员的知识是企业非常重要的财富。如何组织、管理好这些知识财富成为现代管理重要的研究课题^[3,4]。

传统的做法是建立一个数据库系统来管理企业雇员的能力。然而, 基于关键词的查询有时是不能满足人们的要求的。比如查找一个懂得“数据库”的雇员来从事一个项目的开发工作, 假设在建立的能力查询系统中找不到一个人懂得“数据库”, 系统只能返回空记录给用户。而事实上存在某些雇员懂得“关系数据库”或“面向对象数据库”, 按常识, 这些人应该懂得数据库的基本理论, 他们可作为候选人提交给用户。但传统的数据库做不到这一点。当然我们可以通过进一步输入这些关键词达到这一目的, 但我们不能保证用户了解“数据库”所有的子领域。而借助于本体论, 可以实现这一目

作者简介: 廖明宏(1966~), 男, 博士, 主研方向: 人工智能、并行、分布式处理

收稿日期: 1999-08-07

的。假设我们已经为企业建立一个本体论，其中某一部分由图1所示。

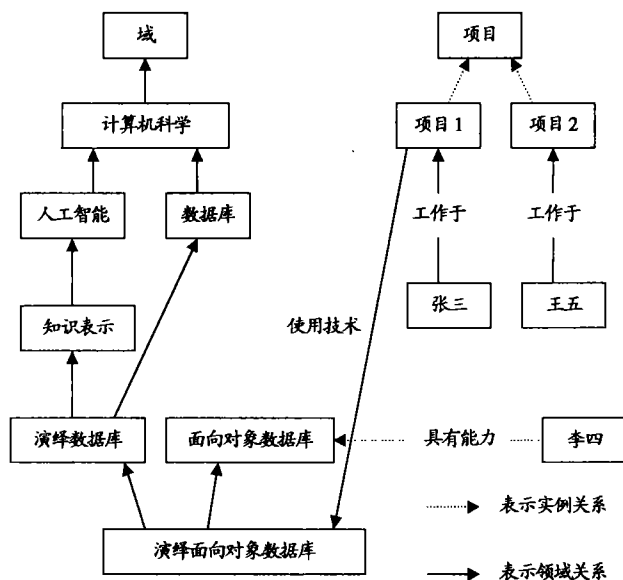


图1 一部分本体论

现在再回到查询懂得数据库的雇员。假设不能得到结果，那么由图1可以看到，李四具有“面向对象数据库”的能力，按常理，他也同样具备“数据库”的能力。一般认为，具备某领域的子领域能力的雇员，同样具备该领域的知识。

我们把问题进一步推广到项目因素。现要查询具备“演绎面向对象数据库”能力的雇员，同样没有结果。然而从图1可以知道张三在干项目1，而项目1应用了“演绎面向对象数据库”技术。因此张三可作为候选人输出。一般地，基于本体论的启发式表达式可定义为 n 个函数的组合运算：

$$f_1 f_2 \dots f_n \quad (1)$$

其中， $f_i(\lambda)^y$

λ 表示一条有向边 e 或其反向(写成 e^{-1})；

y 表示偏序闭包，它表示路径长度，可取下面任一种值： $n, n \dots m, \geq n, *(表示 \geq 0)$ 或 $+(表示 \geq 1)$ 。

公式(1)说明它以一有向图的结点集作为输入，对每个结点，沿着公式所说明的边以从右到左的顺序进行计算，每步产生的中间结果将作为下一步的输入。

例如上述例子可用启发式公式表示为：

- (1) $(has\ Competence^{-1})^1$;
- (2) $(has\ Competence^{-1})^1 (is\ Sub\ Field\ OF^{-1})^*$;
- (3) $(work\ In^{-1})^1 (use\ Technology^{-1})^1$

说明：首先沿 $has\ Competence$ 边的反向直接查找指向某一概念的人；其次查找指向该概念所有子概念的人。启发式查询算法可描述如下：

算法1：启发式信息检索算法

输入：查询关键词集 ks ；启发式公式 $f_1 f_2 \dots f_n$ ；

输出：满足启发式条件的所有记录；

方法：

begin

- (1) 初始化堆栈 $stack$ 为空；
- (2) for $i=1$ to n do

push ($stack, f_i$); // f_i 进栈;

(3) temp= ks ;

(4) while $stack$ 非空 do // 从右向左计算 f_i ;

begin

(5) $f = pop(stack)$;

(6) temp= $f(temp)$;

end;

(7) return temp;

end;

定理1：算法1的时间复杂性为 $O(en)$ ，其中 n 为启发式公式中的函数个数， e 为有向图的边的条数。

证明：算法第(6)步最费时间，由于 $f_i = (\lambda)r$ ，当 r 为 $*$ 时(即求所有子树)，它要求遍历一棵子树，若采用先深遍历算法，其时间复杂性为 $O(e)$ 。算法第(4)步循环 n 次，故整个算法时间复杂性为 $O(en)$ 。

3 一个实验原型系统

为验证上述介绍的启发式信息检索思想，我们基于本体论设计并实现了一个企业雇员能力检索系统(CRS)。该系统基于客户/服务器模型，采用Java编程；因此其客户端可在因特网上用Netscape或Internet Explorer浏览器调用，而服务器部分安装在Sun工作站上。CRS的体系结构见图2。

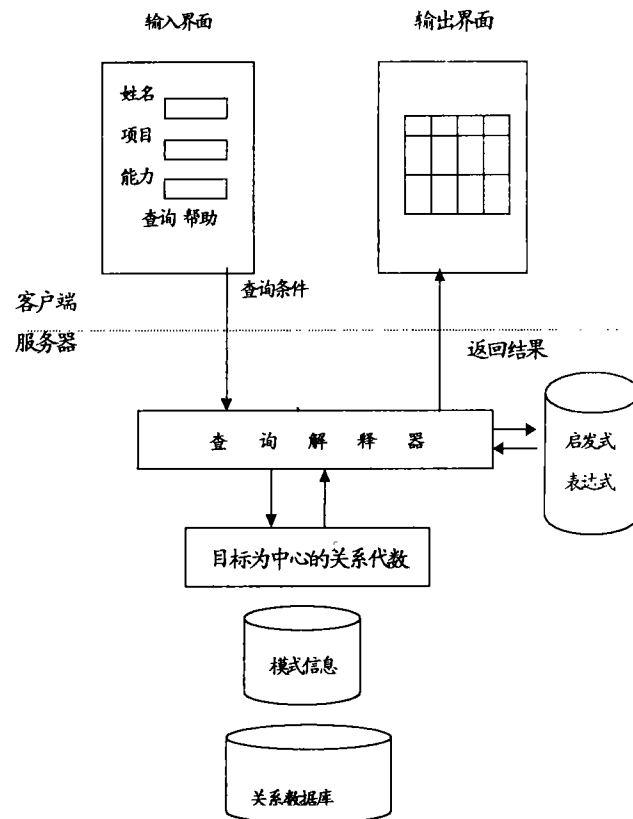


图2 CRS体系结构

在输入部分，用户可查询某一能力，某个人具备的能力或某一项目组成员的能力。其中每个输入域也可以是一个布尔表达式。由客户端提交的查询条件以一定的传输协议送到服务器上。服务器的查询解释根据查询条件和启发式表达式形成关系代数的操作序列，完成相应的操作，并将查询结果

